

MST121 Chapter D4



The Open
University

A first level
interdisciplinary
course

Using
Mathematics

CHAPTER

D4

BLOCK D

MODELLING UNCERTAINTY

Comparing



The Open
University

A first level
interdisciplinary
course

Using **Mathematics**

CHAPTER

D4

BLOCK D

MODELLING UNCERTAINTY

Comparing

Prepared by the course team

About this course

This course, MST121 *Using Mathematics*, and the courses MUI20 *Open Mathematics* and MS221 *Exploring Mathematics* provide a flexible means of entry to university-level mathematics. Further details may be obtained from the address below.

MST121 uses the software program Mathcad (MathSoft, Inc.) and other software to investigate mathematical and statistical concepts and as a tool in problem solving. This software is provided as part of the course, and its use is covered in the associated Computer Book.



The Open University, Walton Hall, Milton Keynes MK7 6AA.

First published 1997. Reprinted 1998, 1999, 2000, 2001.

Copyright © 1997 The Open University

All rights reserved; no part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, or otherwise without either the prior written permission of the Publishers or a licence permitting restricted copying issued by the Copyright Licensing Agency, 90 Tottenham Court Road, London W1P 0LP. This publication may not be lent, resold, hired out or otherwise disposed of by way of trade in any form of binding or cover other than that in which it is published, without the prior consent of the Publishers.

Edited, designed and typeset by The Open University using the Open University T_EX System.

Printed in the United Kingdom by The Burlington Press, Foxton, Cambridge CB2 6SW.

ISBN 0 7492 7895 1

This text forms part of an Open University First Level Course. If you would like a copy of *Studying with The Open University*, please write to the Course Enquiries Data Service, PO Box 625, Dane Road, Milton Keynes MK1 1TY. If you have not already enrolled on the Course and would like to buy this or other Open University material, please write to Open University Educational Enterprises Ltd, 12 Cofferidge Close, Stony Stratford, Milton Keynes MK11 1BY, United Kingdom.

Contents

Study guide	4
Introduction	5
1 Memory and age	6
2 Exploring the data	14
3 How to tell a female meadow pipit from a male	15
4 Testing for a difference	29
Summary of Chapter D4	30
Learning outcomes	30
Solutions to Activities	31
Solutions to Exercises	34

Study guide

You should schedule two study sessions for your work on this chapter, of which one or both will use the computer, depending on the study pattern that you adopt. These sessions may be longer than average. The study pattern which we recommend is as follows.

Study session 1: Sections 1 and 2. You will need access to your computer, together with the statistics software and Computer Book D for Section 2.

Study session 2: Sections 3 and 4. You will need access to your computer, together with the statistics software and Computer Book D for Section 4.

An alternative study pattern, which includes all the computer work in one session, is as follows.

Alternative Study session 1: Sections 1 and 3.

Alternative Study session 2: Sections 2 and 4 (computer).

If you choose this study pattern, then your first session will be a long one.

If you are not very familiar with boxplots and their interpretation, then, whichever study pattern you adopt, you may find that your first session is a long one. If you encountered boxplots for the first time in the preparatory materials, and you feel that you need to consolidate your work on these before beginning this chapter (or during your study of Section 1), then you may wish to make use of the sections on the median and quartiles and on boxplots which are included in *StatsAid*, a short teaching module included with the statistics software. Instructions for using StatsAid are included in an appendix to Computer Book D; whether or not you make use of this software is entirely up to you.

Chapter D5 is designed to be studied in three sessions, each of which should be much shorter than either of the two study sessions for Chapter D4. So, if all the ideas in Chapter D4 are new to you, then you may prefer to study it in three sessions and then study Chapter D5 in two sessions. If you choose to do this, then we suggest you adopt one of the study patterns below.

Study session 1: Sections 1 and 2 (computer),

Study session 2: Section 3,

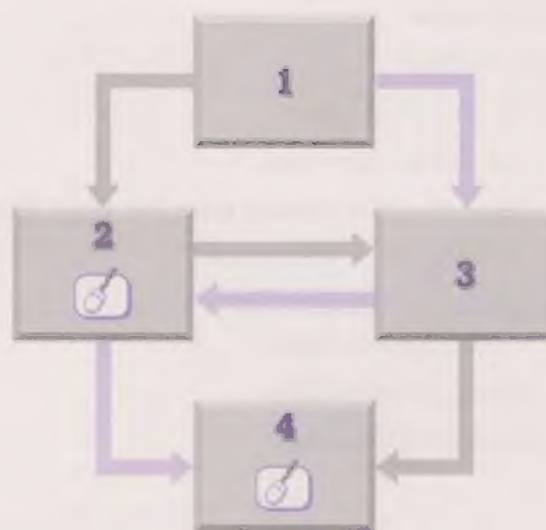
Study session 3: Section 4 (computer);

or

Study session 1: Section 1,

Study session 2: Section 3,

Study session 3: Sections 2 and 4 (computer).



Introduction

Men earn more than women. Non-smokers live longer than smokers. Boys weigh more than girls at birth. Girls are better at reading than boys. The old have poorer memories than the young. How often have you heard comparative statements such as these? And how could you investigate whether there is any truth in them?

In one sense, all the statements are clearly false: it is not true that *all* men earn more than *all* women, for instance, or that *all* non-smokers live longer than *all* smokers. At best, these statements are about averages; for example, the average weight at birth of boys is greater than that of girls.

All the statements above involve a comparison between two groups. To investigate a statement such as ‘boys weigh more than girls at birth’, we would need to obtain the birth weights of some boys and girls; that is, we would need two samples of birth weights. We could then compare them to see if there is evidence of a difference between the birth weights of boys and girls. In this chapter, we look at ways of comparing two samples of data.

In Section 1, an experiment into spatial memory in the young and the elderly is described, and boxplots are used to make a visual comparison of the results obtained for the two groups. Then, in Section 2, the use of OUSats to produce boxplots is discussed. However, boxplots provide only a quick visual comparison of two samples of data. And even if the values in one sample seem to be generally higher than the values in the other, it is possible that this difference is not reflected in the populations from which the samples were drawn: any apparent difference might be due to sampling variation.

But how likely is this to be the case? How much must the samples differ before we can be fairly confident that the populations from which the samples were drawn differ too? To be specific, how large must the difference between the mean birth weight of a sample of boys and the mean birth weight of a sample of girls be before we can conclude that, on average, the mean birth weights of boys and girls are not the same?

This is the type of question that is addressed in Section 3. A procedure for comparing the means of two samples is described: its purpose is to decide whether or not the difference between the *sample* means is large enough to conclude that the *population* means are different. This procedure is an example of a *hypothesis test: the two-sample z-test*. The use of OUSats to perform this test is described in Section 4.

The learning skills theme for this chapter is the same as for Chapter D3: ‘distinguishing different ideas with similar names’.

1 Memory and age

It is commonly believed that as you get older, your memory deteriorates. Is this belief justified? What aspect of memory is referred to here? For instance, many elderly people remember vividly events from their youth, even when the events of last week are forgotten. So we need to distinguish between long-term and short-term memory, as well as between different types of memory – memory of events, of people, of numbers, of words or pictures, and so on. There are many different aspects to memory, so any study of memory and age must be clear about the particular aspect or type of memory that is being investigated.

In the early 1990s, as part of a much broader study, researchers in the Department of Psychology at the University of Sheffield carried out an investigation into spatial memory in the young and the elderly. They were interested in whether there was any difference in the ability of the young and the elderly to remember the positions of objects in space. In one experiment, two groups of people tackled a memory test. Those in one group were aged between 18 and 25 years and those in the other group were over 65 years old; the two groups had similar academic backgrounds. Eighteen everyday objects (a toy car, a thimble, a key, ...) were placed randomly on a 10 by 10 square grid. Each person was asked to study the positions of the objects. When a person indicated that they had looked for long enough, the objects were removed. They were then asked to replace the objects in exactly the same positions.

The data were collected between September 1989 and August 1992 as part of the Ph.D. project of Jennifer Day.

Activity 1.1 Measuring accuracy

Spend a few minutes thinking about how you might measure the accuracy of recall of the participants. Write down your ideas.

Comment

Three possible methods for measuring accuracy of recall are as follows.

- (a) A participant scores 1 for each object that is returned to its original position, and 0 for each object incorrectly placed.
- (b) For each object, a participant scores 1 if it is returned to the correct row, and 1 if it is returned to the correct column (giving a score of 2 for an object that is replaced in its original position).
- (c) For each object, the score awarded is equal to the 'distance' from its original position to its remembered position. One way of measuring this distance is the so-called 'city block score', illustrated in Figure 1.1.

Suppose that an object was originally placed on square *A* and is replaced on square *B*. To get from *A* to *B* involves moving 3 squares horizontally and 2 squares vertically. The city block score for this object is $3 + 2 = 5$. In general, the city block score for an object is the sum of the number of squares horizontally and the number of squares vertically from its original position to its remembered position.

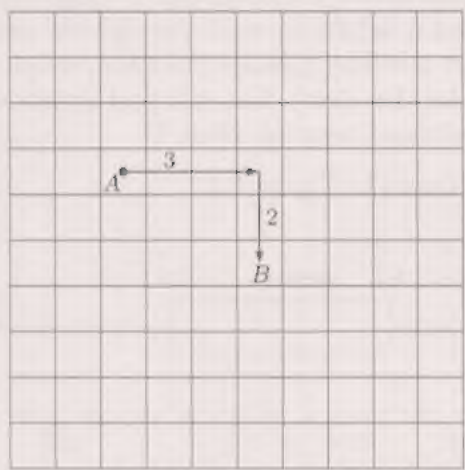


Figure 1.1 Finding a ‘city block score’; $3 + 2 = 5$

You may well have suggested other methods. For instance, you may have suggested using the ‘conventional’ distance: the conventional distance between A and B is $\sqrt{3^2 + 2^2} = \sqrt{13}$. There are many possibilities. Notice that for both the city block score and the conventional distance, a low score corresponds to a good performance on the test.

Activity 1.2 Advantages and disadvantages

What are the advantages and disadvantages of each of the three methods just described, and of any method you suggested for measuring accuracy of recall?

Comment

A solution is given on page 31.

The researchers at the University of Sheffield obtained city block scores for 13 young people and 14 elderly people; the data are given in Table 1.1. Remember that a low score indicates a good performance on the test.

Table 1.1 City block scores

Young	14	29	16	22	11	4	36	6	20	7	12	5	6	
Elderly	17	15	21	34	35	26	32	36	23	42	29	22	13	43

The researchers were interested in whether there is a difference between the ability of the young and of the elderly to remember the positions of the objects. We can investigate this by comparing their city block scores.

In general, it is difficult to draw conclusions simply by inspecting lists of numbers. A useful first step in comparing two sets of data is to make a visual comparison. A diagram which is particularly useful for this purpose is the **boxplot**. If you have studied MU120, then you will be familiar with boxplots and how to interpret them. They are also discussed in the preparatory materials for this course, so only a brief review is included here.

If you are not confident about how to obtain a boxplot to represent a data set, then you may find it helpful to work through the section on boxplots contained in the short teaching package StatsAid, which is provided as part of the software for this block. You will find details of how to use this package in an appendix to Computer Book D.

A typical boxplot is shown in Figure 1.2.

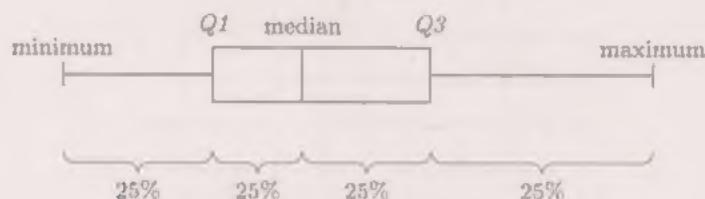


Figure 1.2 A typical boxplot

A boxplot consists of a rectangular box stretching from the *lower quartile* $Q1$ to the *upper quartile* $Q3$, and two whiskers stretching from the ends of the box to the extremes – the minimum and maximum values in the data set. A vertical line is drawn through the box at the median. Roughly speaking, the four parts of the boxplot – the two sections of the box and the two whiskers – each cover approximately 25% of the values in the data set: the lower quartile $Q1$, the median and the upper quartile $Q3$ divide the values in the data set into four subsets, each of which contains approximately 25% of the values. As a reminder (and for your convenience), the definitions of the median and the lower and upper quartiles are given in the box below.

The median

The **median** is essentially the middle value (that is, the middle value when the values are placed in order of size) of a batch of data. It is found by the following procedure.

- ◇ First sort the values into ascending order (if necessary); that is, smallest first, then second smallest, ..., with the largest last.
- ◇ If the batch size is odd, then the median is the middle value in the list.
- ◇ If the batch size is even, then the median is the average (mean) of the two middle values.

The quartiles

Roughly 25% of the values in a batch of data lie below the lower quartile and roughly 25% of the values lie above the upper quartile. The quartiles are defined as follows.

The **lower quartile**, which is denoted $Q1$, is the median of the lower half of the batch – that is, those values to the left of the median when the values in the batch are written in ascending order.

The **upper quartile**, which is denoted $Q3$, is the median of the upper half of the batch.

It is not always possible for *exactly* 25% of the values to lie above the upper quartile; when the batch size n is odd, for instance, $\frac{1}{4}n$ is not a whole number.

Example 1.1 City block scores of the young group

The city block scores of the group of 13 young people are written below in ascending order.

4 5 6 6 7 11 12 14 16 20 22 29 36

Since there is an odd number of values, the median is the middle value, that is, 12. The lower quartile is the median of the values to the left of the median, and the upper quartile is the median of the values to the right of the median. This is illustrated in Figure 1.3.

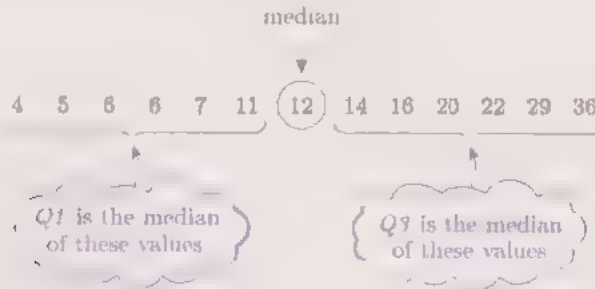


Figure 1.3 Finding the quartiles

So the lower quartile is

$$Q1 = \frac{1}{2}(6 + 6) = 6,$$

and the upper quartile is

$$Q3 = \frac{1}{2}(20 + 22) = 21.$$

Activity 1.3 City block scores of the elderly group

Find the median, the lower quartile and the upper quartile for the city block scores of the group of 14 elderly people.

Comment

The solution is given on page 31.

The results from Example 1.1 and Activity 1.3 were used to produce the boxplots shown in Figure 1.4. Notice that the five key values – the minimum, the lower quartile, the median, the upper quartile and the maximum – have been written on the boxplots. You should always include these values on rough sketches of boxplots, and on accurate boxplots if a scale is not included. So it is a good idea to include them as a matter of course.

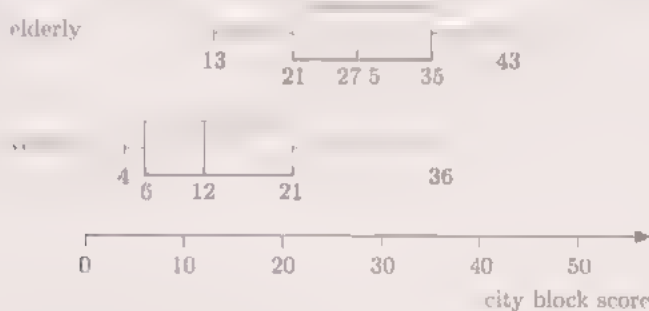


Figure 1.4 Boxplots showing the city block scores

One measure of the difference between the scores of the two groups is given by the location of the boxplots on the city block scale. From the boxplots, it is clear that the scores of the elderly people are generally higher than the scores of the young people, indicating that, on average, the elderly people performed less well on the test than did the young people. All the five key values on a boxplot – the minimum, the lower quartile, the median, the upper quartile and the maximum – are higher for the elderly group than for the young group. In particular, notice that the minimum score for the elderly group (13) is higher than the median score for the young group (12), so more than half of the young people performed better on the memory test than all of the elderly people.

Although the scores of the elderly people are generally higher than the scores of the young people, the boxplots show that the spread of the scores for the elderly group is similar to the spread of the scores for the young group. This can be demonstrated by calculating either of the two measures of spread which may be obtained directly from a boxplot – the *range* and the *interquartile range*.

The **range** is the difference between the maximum and minimum values:

$$\text{range} = \text{maximum} - \text{minimum}.$$

The **interquartile range** is the difference between the upper quartile and the lower quartile:

$$\text{interquartile range} = Q3 - Q1.$$

So, for a boxplot, the length of the box is equal to the interquartile range and the distance from the end of one whisker to the end of the other gives the range. This is illustrated in Figure 1.5.

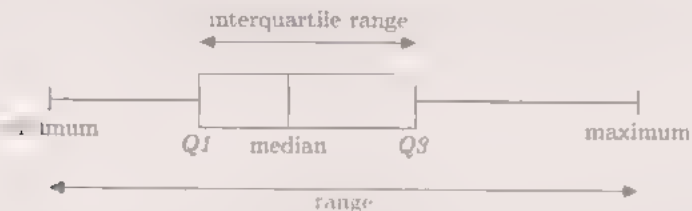


Figure 1.5 Measures of spread

Look again at the boxplots of the city block scores in Figure 1.4. You can see that the boxes are roughly the same length, and the lengths of the boxplots are approximately equal, so the interquartile range and the range are similar for the young group and the elderly group.

Activity 1.4 Measures of spread

Calculate the range and the interquartile range of the city block scores for the group of young people and for the group of elderly people. Do the values confirm that the range and interquartile range are roughly the same for the two groups?

Comment

The solution is given on page 31.

The conclusion we drew from the boxplots of the city block scores is that there seems to be a difference between the ability of the young and the elderly to remember the positions of the objects on the grid. In fact, the young people seem to do better on the test – their city block scores are generally lower. However, when the experiment was carried out, the participants were allowed to study the positions of the objects for as long as they wished. It is possible that the longer you spend studying the positions of the objects, the better you will remember them. So did the two groups spend similar lengths of time studying the positions of the objects? If the young people spent longer than the elderly, then this by itself could explain why their scores were lower. In the next activity, you are asked to compare the times the two groups spent studying the positions of the objects.

Activity 1.5 Comparing memorisation times

Table 1.2 shows the times spent studying the positions of the objects by the 13 young people and the 14 elderly people.

Table 1.2 Memorisation times in seconds

Young	90	90	100	55	145	130	55	85	95	140	125	70	105	
Elderly	75	90	40	40	25	30	55	45	35	55	35	100	45	40

- For each group, sort the times into ascending order, and find the median, the lower quartile and the upper quartile.
- Draw boxplots for the memorisation times of the two groups, using a common axis.
- What do the boxplots tell you about the times spent by the two groups memorising the positions of the objects?

Comment

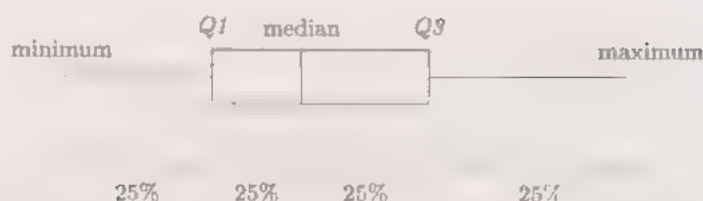
The solution is given on page 31.

This section has included a brief review of the use of boxplots to compare visually two sets of data. Two exercises are provided for you to carry out if you need further practice at working out medians and quartiles or drawing boxplots. In the next section, you will learn how to obtain boxplots using OUStats; you will have further opportunities there to compare two sets of data by interpreting boxplots drawn on a common axis. However, even when boxplots suggest that there is a difference between two sets of data – such as, for example, between the city block scores of the young and the elderly – you need to consider whether the apparent difference could simply be the result of sampling variation.

In Chapter D3, you saw that there can be considerable variation between different samples drawn from the same population. But how different must two samples be before we can be confident that they do not come from the same population? In the context of the memory tests, how different must the scores of the young group and the elderly group be before we can be confident that young people in general are better than the elderly at memorising the positions of objects in space? This is the problem that we shall address in Section 3, where a test is introduced for investigating the difference between two population means given a sample of data from each population.

Summary of Section 1

A typical boxplot is shown below.



The five values marked on a boxplot are the minimum, the lower quartile $Q1$, the median, the upper quartile $Q3$ and the maximum. The lower quartile, the median and the upper quartile divide a batch of data into four parts, each of which contains approximately 25% of the values in the batch.

If the values in a batch of data are written in ascending order and the batch size is odd, then the median is the middle value in the list; if the batch size is even, then the median is the mean of the two middle values.

The lower quartile is the median of the lower half of the batch – that is, those values to the left of the median when the values in the batch are written in ascending order.

The upper quartile is the median of the upper half of the batch.

The range is the difference between the maximum and minimum values:

$$\text{range} = \text{maximum} - \text{minimum}.$$

The interquartile range is the difference between the upper quartile and the lower quartile:

$$\text{interquartile range} = Q3 - Q1.$$

Exercises for Section 1

Exercise 1.1

The table below gives the gross weekly earnings, including overtime (in pounds), of 19 police officers (sergeants and constables) in 1995.

Table 1.3

Women	360	405	315	390	495	430	330	455	365	
Men	455	340	530	440	425	485	355	420	550	400

- Find the median, the lower quartile and the upper quartile for the earnings of the 9 women and, separately, for the earnings of the 10 men.
- Draw boxplots for the gross weekly earnings of the men and of the women.
- What do the boxplots tell you about the relative earnings in 1995 of male and female sergeants and constables?
- Calculate the range and the interquartile range of the women's earnings and of the men's earnings. Is the spread of the women's earnings greater or less than the spread of the men's earnings?

Exercise 1.2

The table below gives the gross hourly earnings, including overtime (in pence), of 19 chefs and cooks in 1995.

Table 1.4

Women	445	325	570	380	315	485	295	370				
Men	550	505	430	620	640	830	360	340	750	405	490	

- Find the median, the lower quartile and the upper quartile for the earnings of the 8 women and, separately, for the earnings of the 11 men.
- Draw boxplots for the gross hourly earnings of the men and of the women.
- What do the boxplots tell you about the relative earnings in 1995 of male and female chefs and cooks?
- Calculate the range and the interquartile range of the women's earnings and of the men's earnings. Is the spread of the women's earnings greater or less than the spread of the men's earnings?

2 Exploring the data



To study this section, you will need access to your computer and the statistics software.

In this section, the use of OUStats to produce boxplots is illustrated for the data on city block scores in Table 1.1. You will be invited to explore the data further to see whether there is a relationship between the time spent memorising the positions of the objects and the score obtained on the test.



Refer to Computer Book D for the work in this section.

Summary of Section 2

In this section, OUStats has been used to investigate further the data on city block scores and memorisation times for young and elderly people. The use of OUStats to produce boxplots has been described.

3 *How to tell a female meadow pipit from a male*

In many studies of bird behaviour, it is important to be able to determine whether a particular bird is male or female. For some species, the sex of a bird is obvious to the observer. For example, whereas adult male blackbirds are truly black, the adult females are brown. For other species, the differences in plumage are more subtle, and it is only when a close examination can be made that an ornithologist is able to determine the sex. For example, fieldfares cannot be sexed by the distant observer, but once a fieldfare is held in the hand, the pattern of the crown feathers provides an effective criterion for determining sex.

In the breeding season, the sex of many species of birds can be determined by examining the shape of the underparts of their bodies. In addition, in many species, just before incubation starts, the females shed the downy feathers on their underparts, producing what is known as an incubation patch.

However, there are many species of bird, such as robins and meadow pipits, where, outside the breeding season, there is no visible difference in plumage between males and females, even when they are examined by hand. So how can an ornithologist tell whether a particular bird is male or female?

One way that this problem has been tackled is by taking measurements *during* the breeding season for various features, such as wing length and weight, for birds of known sex. The data collected are then analysed for any differences between these features for males and females. If a marked difference is found in some measurement, then this could lead to a way of determining sex outside the breeding season.

In one study, the wing lengths of 31 male and 27 female meadow pipits were measured to the nearest millimetre during the breeding season. The data are given in Table 3.1.

Table 3.1 Wing lengths of meadow pipits in millimetres

Males	81	84	79	84	78	81	82	83	85	80	81
	81	79	79	80	81	82	83	81	84	81	
	83	82	83	82	82	79	79	83	83	81	
Females	77	77	80	76	80	78	77	80	77	80	
	79	75	77	79	77	76	75	75	75	80	
	76	81	82	75	77	78	74				

Boxplots for the wing lengths of these meadow pipits are shown in Figure 3.1.

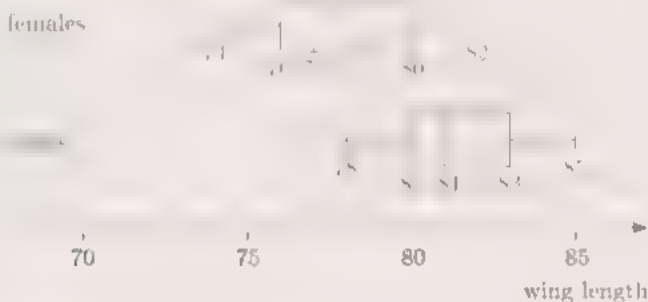


Figure 3.1 Wing lengths of meadow pipits

Activity 3.1 Interpreting the boxplots

What do the boxplots tell you about the wing lengths of male and female meadow pipits?

Comment

The wing lengths of the males were in general greater than the wing lengths of the females, although there was considerable overlap. For instance, the wing lengths of approximately 25% of the males were less than 80 mm and the wing lengths of approximately 25% of the females were greater than 80 mm.

'Population' has its everyday meaning here.

From the boxplots, it looks as though, on average, the wing lengths of males are greater than the wing lengths of females. However, we cannot be certain that this is the case: it is possible that the mean wing length of the population of male meadow pipits is equal to the mean wing length of the population of female meadow pipits, and that the difference in mean wing lengths observed in these samples is simply due to sampling variation. But is this likely? How much must the wing lengths in the two samples differ before we can be confident that there really is a difference between the average wing lengths of male and female meadow pipits? More specifically, how different must the sample means be before we can be confident that the population means are different?

This sort of question can be answered by carrying out a statistical procedure called a *hypothesis test*. In this section, a hypothesis test called the *two-sample z-test* is described. The problem of deciding whether or not we can be confident that the mean wing lengths of male and female meadow pipits differ will be used to illustrate the main features of the test.

There are three main stages involved in carrying out a hypothesis test: setting up hypotheses to be tested, calculating a number called the *test statistic*, and drawing conclusions. We shall discuss each of these in turn, using the wing lengths of male and female meadow pipits to illustrate the ideas. We would like you to get an idea of what is involved in carrying out a hypothesis test from beginning to end, so we shall not interrupt this discussion with activities based on other investigations: most of the activities are at the end of this section. However, try to read the next few pages actively: make a note of new terminology and notation, and make sure you understand the explanations. You may wish to re-read some of the discussion when you come to carrying out a hypothesis test yourself.

A word of warning is appropriate here. Do not expect to master the idea of a hypothesis test at first reading if it is new to you: it is a major idea, and it may take you some time to grasp it fully. In this section, our aim is simply to introduce you to hypothesis testing so that when you meet other hypothesis tests you will understand the principles and the terminology involved. The best way of assimilating the ideas is by carrying out tests, so you should find working through the examples and activities in this section and the next very helpful. If you do find the ideas difficult to grasp, then try re-reading this section *after* you have worked through the activities.

Hypotheses

A hypothesis test begins with some sort of hypothesis about the population or populations of interest. In this case, there are two populations – the wing lengths of male and female meadow pipits. We want to know whether or not the mean wing length for male meadow pipits is equal to the mean wing length for females. So our hypothesis is that

‘Population’ has its technical meaning here

the mean wing lengths of males and females are equal,

that is, they do not differ. This hypothesis is called the **null hypothesis** for the test, and is usually denoted by H_0 : the word ‘null’ is used because it is a hypothesis of ‘no’ difference. So the null hypothesis for the test can be written as follows.

H_0 : The mean wing length of male meadow pipits is equal to the mean wing length of female meadow pipits.

At the end of a hypothesis test, we either accept the null hypothesis or we reject it in favour of what is called the **alternative hypothesis**. In this example, the alternative hypothesis is that the mean wing length of male meadow pipits is not equal to the mean wing length of female meadow pipits. The alternative hypothesis is usually denoted by H_1 , so we can write it as follows.

H_1 : The mean wing length of male meadow pipits is not equal to the mean wing length of female meadow pipits.

Every hypothesis test should begin with a statement of two hypotheses: the null hypothesis H_0 and the alternative hypothesis H_1 . From now on we shall use this standard terminology: we shall always begin a hypothesis test by stating clearly the null hypothesis H_0 and the alternative hypothesis H_1 .

It is convenient to introduce symbols for the population means and standard deviations. This will enable us to express the null and alternative hypotheses more concisely. And we shall need these symbols in order to explain the details of the test. We shall denote the means of the two

populations μ_M and μ_F and the standard deviations of the populations by σ_M and σ_F . So, for instance, μ_M is the mean wing length of the population of male meadow pipits, and μ_F is the mean wing length of the population of female meadow pipits. Using μ_M and μ_F , we can write the hypotheses in the following concise form:

Subscript M for male, F for female.

$$H_0: \mu_M = \mu_F,$$

$$H_1: \mu_M \neq \mu_F.$$

Before considering the second stage of a hypothesis test, pause for a moment to note the dual use of the word ‘population’ in the preceding text. It has been used in its everyday sense – the population of female meadow pipits – and in its technical sense – the population of wing lengths of female meadow pipits. This dual use is common, and you should not be concerned about it: it should not lead to confusion as it is usually clear when the word ‘population’ is being used in its technical sense.

The test statistic

If the population means are equal, that is, if the null hypothesis is true, then we would not expect the sample means to differ greatly. If the sample means do differ greatly, then this would be evidence against the population means being equal, that is, against the null hypothesis and in favour of the alternative hypothesis that the population means are not equal. So we need to look at the difference between the sample means, $\bar{x}_M - \bar{x}_F$, and assess whether this difference is 'large enough' to reject the null hypothesis. This will involve using the results for sampling distributions that you met in Chapter D3. Although you will not be expected to reproduce the details of the arguments presented in the next few pages, do try to follow them. If you understand how the final result is derived, then this will give you a greater appreciation of how a hypothesis test works.

First, we shall summarise the data in the samples. The sample means, \bar{x}_M and \bar{x}_F , the sample standard deviations, s_M and s_F , and the sample sizes, n_M and n_F , are given in Table 3.2.

Table 3.2 Wing lengths of meadow pipits (in millimetres)

	Sample size	Sample mean	Sample standard deviation
Males	$n_M = 31$	$\bar{x}_M = 81.5$	$s_M = 1.79$
Females	$n_F = 27$	$\bar{x}_F = 77.5$	$s_F = 2.15$

Before we can assess whether the difference between the sample means $\bar{x}_M - \bar{x}_F$ is 'large enough' for us to reject the null hypothesis, we need a result concerning the *sampling distribution of the difference between two sample means*.

If a sample is drawn from each population, and the sample means x_M and x_F are found, then the difference $\bar{x}_M - x_F$ can be calculated. For different pairs of samples, the differences $\bar{x}_M - x_F$ will vary. Imagine that this difference could be calculated for all possible pairs of samples. Then the distribution of the differences is the **sampling distribution of the difference between two sample means**.

You already know that, by the Central Limit Theorem, if the sample size is fairly large (at least 25), then the sampling distribution of the mean is approximately a normal distribution. Moreover, the sampling distribution of the mean for samples of n_M wing lengths of males has mean μ_M and standard deviation $\sigma_M/\sqrt{n_M}$, and the sampling distribution of the mean for samples of n_F wing lengths of females has mean μ_F and standard deviation $\sigma_F/\sqrt{n_F}$.

The result that we need here depends on the above results and also on a result which states that the distribution of the difference between two (independent) random variables which are each normally distributed is also a normal distribution. Thus we have the following result.

Provided that the sample sizes are sufficiently large (at least 25), the sampling distribution of the difference between two sample means is approximately a normal distribution.

For the wing lengths of the meadow pipits, the sample sizes are $n_M = 31$ and $n_F = 27$, so we can assume that the sampling distribution of the difference between two sample means is approximately normal. We also

Here, 'large' means either 'large and negative' or 'large and positive'.

The Central Limit Theorem is discussed in Chapter D3, Subsection 1.3.

need to know the mean and standard deviation of this sampling distribution. The mean of the sampling distribution is equal to the difference between the population means, $\mu_M - \mu_F$. And its standard deviation, which is called the **standard error of the difference between two sample means**, is given by

$$SE = \sqrt{\frac{\sigma_M^2}{n_M} + \frac{\sigma_F^2}{n_F}}.$$

Notice that this formula involves the standard deviations of the two sampling distributions of the mean: $\sigma_M/\sqrt{n_M}$ and $\sigma_F/\sqrt{n_F}$. The derivation of the formula requires results which are beyond the scope of this course. However, if you study statistics further, then you will almost certainly encounter the necessary techniques and results. The sampling distribution is summarised in the box below and illustrated in Figure 3.2.

The full details are given in the Open University course M246, for instance.

The sampling distribution of the difference between two sample means

For samples of sizes n_M, n_F (where n_M, n_F are both at least 25) from populations with means μ_M, μ_F and standard deviations σ_M, σ_F , the sampling distribution of the difference between two sample means is approximately a normal distribution with mean $\mu_M - \mu_F$ and standard deviation given by

$$SE = \sqrt{\frac{\sigma_M^2}{n_M} + \frac{\sigma_F^2}{n_F}}.$$

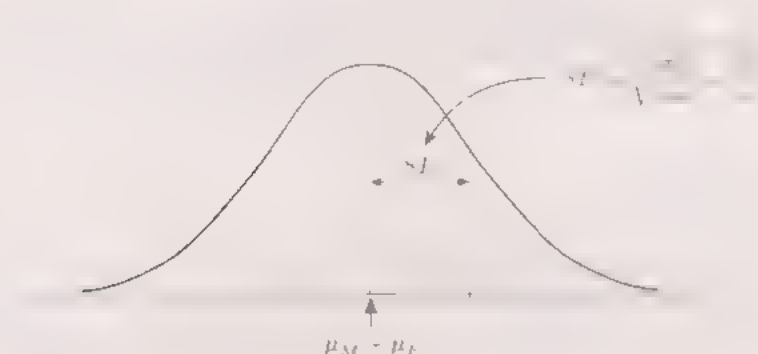


Figure 3.2 The sampling distribution of the difference between two sample means

If the null hypothesis is true, that is, if the population means μ_M and μ_F are equal, then the distribution of the difference between two sample means will have mean 0 (since in that case $\mu_M - \mu_F = 0$). So, if the null hypothesis is true, the sampling distribution is approximately normal with mean 0 and standard deviation given by

$$SE = \sqrt{\frac{\sigma_M^2}{n_M} + \frac{\sigma_F^2}{n_F}}.$$

You know that for *any* normal distribution, 95% of values lie within 1.96 standard deviations of the mean, and 5% of values lie 1.96 or more standard deviations from the mean (see Figure 3.3).

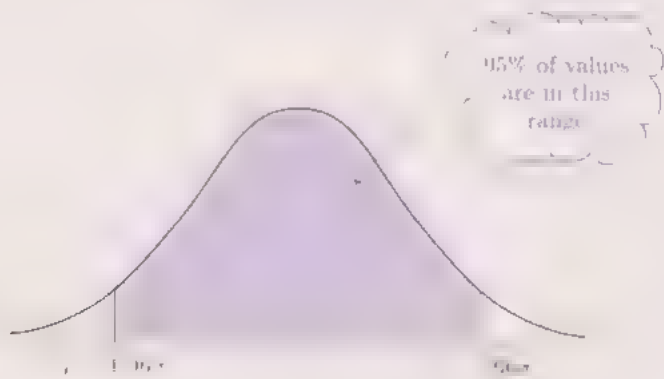


Figure 3.3 A normal distribution

So 95% of differences $\bar{x}_M - \bar{x}_F$ will be within 1.96 standard deviations of the mean $\mu_M - \mu_F$ (see Figure 3.4).

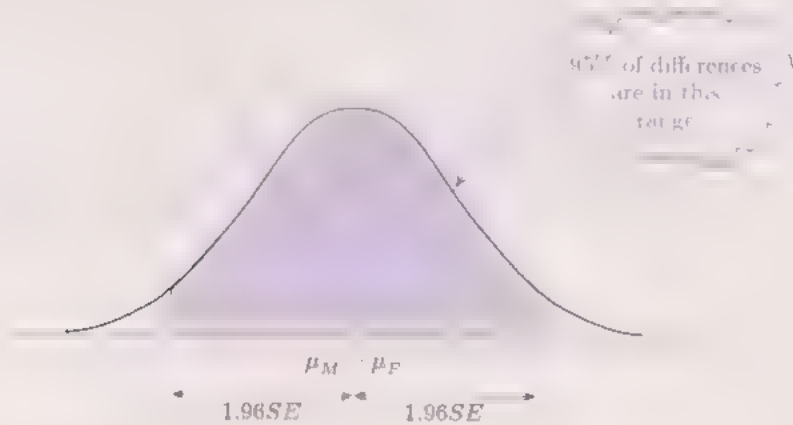


Figure 3.4 Values of $\bar{x}_M - \bar{x}_F$

So if the null hypothesis is true, then 95% of differences $\bar{x}_M - \bar{x}_F$ will be within 1.96 standard deviations of 0 (see Figure 3.5).

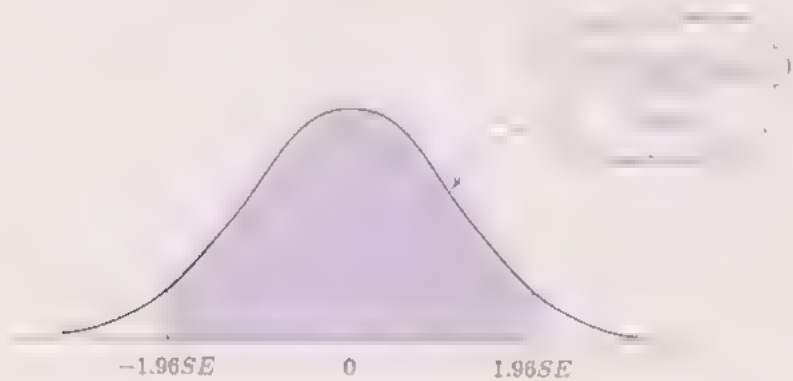


Figure 3.5 Values of $\bar{x}_M - \bar{x}_F$ when the null hypothesis is true

Equivalently, there is only a 5% chance that the sample means will differ by at least 1.96 standard deviations. So a difference as large or larger than 1.96 standard deviations is unlikely to occur if the population means are equal. Hence if a difference of this size does occur, then we might reasonably doubt that the population means are equal. This is the basis of

our hypothesis test. If the difference between the sample means is 'large', then we reject the null hypothesis that the population means are equal.

Essentially, the test involves finding the difference between the sample means; then if this difference is at least 1.96 standard deviations (that is, $1.96SE$), we reject the null hypothesis.

The difference between the sample means is $\bar{x}_M - \bar{x}_F$. We shall reject the null hypothesis if this difference is in either of the shaded areas in Figure 3.6. These areas correspond to 'large' differences.

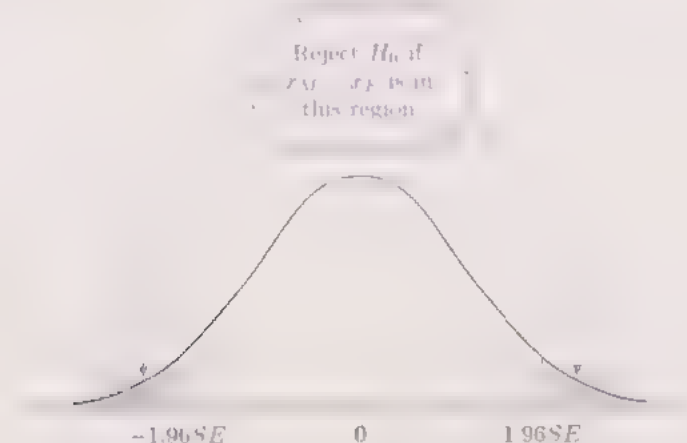


Figure 3.6 Differences for which H_0 is rejected

So we must compare the difference between the sample means $\bar{x}_M - \bar{x}_F$ with the standard error SE ; that is, we must find $(\bar{x}_M - \bar{x}_F)/SE$. Then we shall reject H_0 if

$$\text{either } \frac{\bar{x}_M - \bar{x}_F}{SE} \leq -1.96 \quad \text{or} \quad \frac{\bar{x}_M - \bar{x}_F}{SE} \geq 1.96$$

However, we cannot calculate $(\bar{x}_M - \bar{x}_F)/SE$ because we do not know σ_M and σ_F , the two population standard deviations, and hence we do not know the value of SE . We deal with this problem in exactly the same way as we did in Chapter D3, when calculating confidence intervals: we replace σ_M by s_M , and σ_F by s_F ; that is, we use the sample standard deviations to estimate the (unknown) population standard deviations. This gives us an estimated value for SE , the standard deviation of the sampling distribution of the difference between two sample means, which we shall denote by ESE for convenience:

$$ESE = \sqrt{\frac{s_M^2}{n_M} + \frac{s_F^2}{n_F}}$$

We shall refer to this as the **estimated standard error**. So, instead of $(\bar{x}_M - \bar{x}_F)/SE$, we shall calculate $(\bar{x}_M - \bar{x}_F)/ESE$. This quantity is the **test statistic** for the test; it is usually denoted by z . We shall reject the null hypothesis H_0 if the test statistic z is 'large' or, more precisely, if

$$\text{either } z \leq -1.96 \quad \text{or} \quad z \geq 1.96,$$

where

$$z = \frac{\bar{x}_M - \bar{x}_F}{ESE}.$$

If $-1.96 < z < 1.96$, then we shall not reject the null hypothesis. (Notice that when z is equal to -1.96 or 1.96 , we reject the null hypothesis.)

Activity 3.2 *Calculating the test statistic*

- (a) Use the summary statistics in Table 3.2 for the wing lengths of male and female meadow pipits to calculate the estimated standard error.
- (b) Calculate the test statistic.
- (c) Should the null hypothesis be rejected?

Comment

The solution is given on page 32.

Conclusions

In this example, the test statistic is 'large', that is, at least 1.96 in size. So we reject the null hypothesis H_0 . Although a 'large' test statistic is unlikely if the population means are equal, approximately 5% of differences lead to a 'large' test statistic, so there is a 5% chance that we shall wrongly reject the null hypothesis. We say that the **significance level** of the test is 5%. If a test uses a 5% significance level, then this means that there is a 5% chance that we shall wrongly reject the null hypothesis.

It is possible to use a different significance level. For confidence intervals, using a confidence level other than 95% required a different value from the standard normal distribution in place of 1.96. Similarly, for the two-sample z -test, to use a significance level other than 5%, we need to use a different value obtained from the standard normal distribution in place of 1.96. For instance, if we wanted the chance of wrongly rejecting the null hypothesis to be only 1%, then we would use a 1% significance level. This might be the case, for instance, if we do not want to reject the null hypothesis unless there is very strong evidence that it is false. In this case, since 99% of values in a normal distribution lie within 2.58 standard deviations of the mean, we would replace 1.96 with 2.58. If you study statistics further, then you will learn how to use various significance levels. However, in this course, we shall use only a 5% significance level.

Since we are using a 5% significance level, our conclusion should include this information. So, instead of saying simply 'we reject H_0 ', we should say 'we reject H_0 at the 5% significance level'. And instead of ' H_0 is not rejected', we should say ' H_0 is not rejected at the 5% significance level'.

You might think that once H_0 has been rejected or not, as the case may be, we have completed the test. If so, you would be wrong. It is essential to remember what the hypotheses are that have been tested. You should express your conclusion in terms of these hypotheses. For example, for the wing lengths of meadow pipits, the conclusion could be as follows.

Since $z = 7.63 > 1.96$, we reject the null hypothesis at the 5% significance level in favour of the alternative hypothesis. We conclude that the mean wing length of male meadow pipits is not equal to the mean wing length of female meadow pipits. Moreover, the sample mean is greater for the males than for the females, so this suggests that the mean wing length of males is greater than the mean wing length of females.

Notice the final sentence. Having concluded that there is a difference between the population means, we looked at the data to see which sample mean is the greater. Since the sample mean is greater for the males than

for the females, it seems likely that the population mean is greater for the males than for the females, rather than vice versa. The final sentence of the conclusion above is just a statement of this commonsense deduction. If you reject the null hypothesis that the population means are equal, then it is good practice to look again at the data to see which population mean is likely to be the greater.

The two-sample z-test: a summary

The hypothesis test that has been described is called the **two-sample z-test**. The development of the test used the Central Limit Theorem, so the test can be applied only when the sample sizes are large: both sample sizes should be at least 25. (If either sample size is less than 25, then a different test must be used; we shall not be discussing other hypothesis tests in this course.) The three main stages in carrying out the two-sample z-test have been discussed in quite a lot of detail, in order to explain the ideas behind the test. In practice, the test is straightforward to apply. The procedure is summarised in the following box. The populations are labelled A and B , so that the summary is quite general.

Procedure for the two-sample z-test

Stage 1: Hypotheses

Set up the null and alternative hypotheses:

$$H_0: \mu_A = \mu_B,$$

$$H_1: \mu_A \neq \mu_B,$$

where μ_A and μ_B are the means of populations A and B , respectively.

Stage 2: The test statistic

Calculate the test statistic

$$z = \frac{(\bar{x}_A - \bar{x}_B) - (\mu_A - \mu_B)}{ESE}$$

where

$$ESE = \sqrt{\frac{s_A^2}{n_A} + \frac{s_B^2}{n_B}},$$

\bar{x}_A and \bar{x}_B are the sample means, s_A and s_B are the sample standard deviations, and n_A and n_B are the sizes of the samples from populations A and B , respectively.

Stage 3: Conclusions

- ◇ If $z \leq -1.96$ or $z \geq 1.96$, then H_0 is rejected at the 5% significance level in favour of the alternative hypothesis.
- ◇ If $-1.96 < z < 1.96$, then H_0 is not rejected at the 5% significance level.

The conclusion should be expressed in terms of the hypothesis being tested.

Example 3.1 Wing lengths of fieldfares

During the winters of 1980–89, just over 1000 fieldfares were caught in an orchard in Daresbury, Cheshire. All the birds were sexed and aged as first-years or adults (any bird more than one year old). For some of the birds, one or both of wing length and weight were measured. Wing lengths were measured to the nearest millimetre, and weights to the nearest gram. A summary of the data collected on wing lengths of 594 fieldfares is given in Table 3.3.

Table 3.3 Wing lengths of fieldfares (in millimetres)

	Sample size	Sample mean	Sample standard deviation
Adult males	80	151.9	3.19
Adult females	128	147.5	3.37
First-year males	131	150.0	3.10
First-year females	255	146.1	3.37

The sample sizes are large, so we can use the two-sample z -test. We shall test the hypothesis that there is no difference between the mean wing length of adult male fieldfares and the mean wing length of adult female fieldfares.

Using μ_{AM} for the mean wing length of adult male fieldfares and μ_{AF} for the mean wing length of adult female fieldfares, the null and alternative hypotheses can be written as

$$H_0 : \mu_{AM} = \mu_{AF},$$

$$H_1 : \mu_{AM} \neq \mu_{AF}.$$

The estimated standard deviation of the sampling distribution of the difference between two sample means is

$$ESE = \sqrt{\frac{s_{AM}^2}{n_{AM}} + \frac{s_{AF}^2}{n_{AF}}} = \sqrt{\frac{3.19^2}{80} + \frac{3.37^2}{128}} = 0.464679\dots$$

So the test statistic is

$$= \frac{\bar{x}_{AM} - \bar{x}_{AF}}{ESE} = \frac{151.9 - 147.5}{0.464679\dots} \simeq 9.47$$

Since the test statistic is $z = 9.47 > 1.96$, we reject the null hypothesis at the 5% significance level in favour of the alternative hypothesis. We conclude that the mean wing length of adult male fieldfares is not equal to the mean wing length of adult female fieldfares. The sample mean is greater for the males than for the females, so this suggests that the mean wing length of adult males is greater than the mean wing length of adult females.

Activity 3.3 Wing length and age

Use the data in Table 3.3 and the two-sample z -test to investigate whether there is any difference between the mean wing lengths of:

- (a) first-year male fieldfares and adult male fieldfares;
- (b) first-year female fieldfares and adult female fieldfares.

In each case, be sure to specify the null and alternative hypotheses, calculate the test statistic, and state your conclusion clearly.

Comment

The solution is given on page 32.

Source: David Norman (1995) 'Flock composition and biometrics of Fieldfares *Turdus pilaris* wintering in a Cheshire orchard', *Ringing and Migration*, 16, pp. 1–13.

In this case, the populations consisted of all male and all female fieldfares that wintered in the orchard.

Activity 3.4 Weights of fieldfares

The data on the weights of 664 of the fieldfares caught in the same Cheshire orchard are summarised in Table 3.4.

Table 3.4 Weights of fieldfares in grams

	Sample size	Sample mean	Sample standard deviation
Adult males	93	114.9	10.91
Adult females	139	108.7	9.06
First-year males	144	111.6	8.62
First-year females	288	108.0	9.52

Use the two-sample *z*-test to investigate whether there is any difference between the mean weight of first-year female fieldfares and the mean weight of adult female fieldfares. (In Exercise 3.2, you will be asked to investigate whether there are any differences between the mean weights of the other categories of fieldfares.)

Comment

The solution is given on page 33.

Activity 3.5 Another standard error

In Activities 1.7 and 2.8 of Chapter D3, you were asked to construct a table summarising information and results about several different standard deviations: the population standard deviation, the standard error of the mean and the sample standard deviation. In this section, another standard deviation has been introduced: the standard error of the difference between two sample means. Add this standard error to your table. Also include in your table an entry for the estimated standard error of the difference between two sample means, *ESE*. If there are any notes or results that you wish to add to your summary, then do so now.

Comment

Some comments are given on page 33.

Postscript: First-year meadow pipits

In Chapter D2, data were given on the wing lengths in millimetres of 252 first-year meadow pipits caught at Leadburn in southern Scotland in the autumn of 1991. These data are represented in Figure 3.7.

See Chapter D2, Activity 2.2

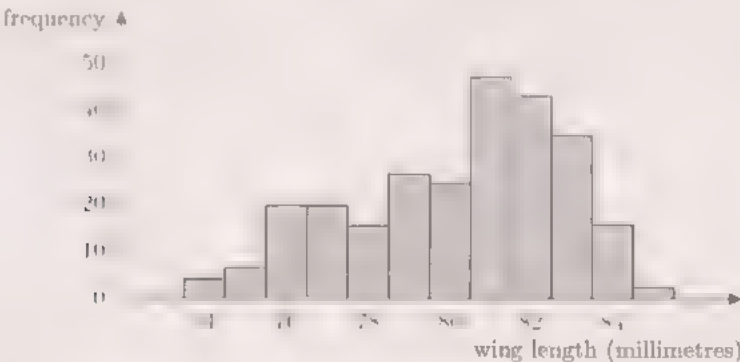


Figure 3.7 Wing lengths of first-year meadow pipits

One of the purposes of the study in which these birds were caught was to estimate the proportions of males and females among first-year meadow pipits. But, as already mentioned, it is not possible to tell the sex of a meadow pipit simply by observation. However, it has been established in a number of studies that the wing lengths of males are, on average, greater than the wing lengths of females, although, as you saw in Figure 3.1, the ranges of wing lengths for the two sexes overlap.

Look again at Figure 3.7. Since the data include measurements on the wing lengths of both male and female meadow pipits, the diagram represents measurements taken from two populations rather than just one. The frequency diagram has a clear peak at about 81 mm, and there is a suggestion of a smaller peak at about 76 or 77 mm. It is possible that these peaks correspond to peaks in the two separate populations. Now that we know that the average wing lengths of male and female meadow pipits differ, we can see that the model suggested in Chapter D2 would be inappropriate: two models are needed, one for males and one for females.

The fact that the wing lengths of males and females differ suggests a practical method for sexing at least some meadow pipits — those with the longest wing lengths are classified as male, those with the shortest wing lengths as female, and those with intermediate wing length are left unclassified. Before the birds can be sexed, a classification rule must be decided. Since wing lengths may vary a little from one bird population to another, the rule is not based on data from other populations. In this study, the data in Figure 3.7 were themselves used to formulate a rule.

For situations where it is known that the measurements collected are from two populations and not just one (the two populations being the wing lengths of males and females in this case), special graphical techniques have been developed for estimating the means and standard deviations of the two populations. For the Leadburn data, these techniques produced estimates for the mean wing lengths of males and females of 81.7 mm and 77.5 mm, respectively. The corresponding estimates for the standard deviations were 1.4 mm and 1.9 mm. These estimates were used to determine a rule for sexing the birds. Birds with wing lengths less than or equal to 79 mm were classified as female and birds with wing lengths greater than or equal to 81 mm were classified as male. The remaining birds — that is, those with wing lengths measured as 80 mm — were left unsexed. (The wing lengths were measured to the nearest millimetre.) Using this rule, roughly 56% of the birds were classified as male, 35% as female, and the rest, 9%, were not sexed. No evidence was found that the techniques used to trap the birds produced any bias in the results: a male was not more likely to be caught than a female. Thus the data provide some evidence that this particular population comprised a greater proportion of male birds than female birds. This suggests that there were more males than females among juvenile meadow pipits in this autumn population.

We shall not describe such techniques here.

Summary of Section 3

In this section, you have been introduced to hypothesis testing; the two-sample z -test has been discussed in some detail. To use this test, a sample of at least 25 measurements is required from each of two populations. The test may be used to investigate whether there is a difference between the means of the populations.

We begin a hypothesis test by specifying appropriate null and alternative hypotheses. For the two-sample z -test, the null hypothesis is that the population means are equal. The alternative hypothesis is that the population means are not equal. That is,

$$H_0: \mu_A = \mu_B,$$

$$H_1: \mu_A \neq \mu_B,$$

where μ_A and μ_B are the means of the two populations A and B .

The next step is to calculate the test statistic. For this test, the test statistic is

$$ESE = z,$$

where

$$ESE = \sqrt{\frac{s_A^2}{n_A} + \frac{s_B^2}{n_B}}$$

\bar{x}_A and \bar{x}_B are the means of the samples from populations A and B , s_A and s_B are the sample standard deviations, and n_A and n_B are the sample sizes.

If the test statistic z satisfies $z \leq -1.96$ or $z \geq 1.96$, then the null hypothesis is rejected at the 5% significance level in favour of the alternative hypothesis

If $-1.96 < z < 1.96$, then the null hypothesis is not rejected at the 5% significance level.

When using a 5% significance level, as has been done in this section, there is a 5% chance that the null hypothesis will be wrongly rejected.

Exercises for Section 3

Exercise 3.1 *Weights of aquatic warblers*

See Chapter D3, Example 2.1 and Activity 2.3.

In Section 2 of Chapter D3, data collected in north-eastern Poland in the early 1990s were used to calculate confidence intervals for the mean weights of male and female aquatic warblers. The confidence intervals did not overlap. It looked as though male aquatic warblers are heavier, on average, than females. The data are summarised in Table 3.5.

Table 3.5 Weights of aquatic warblers in grams

	Sample size	Sample mean	Sample standard deviation
Males	66	12.6	0.73
Females	83	12.1	0.87

Use the two-sample z -test to investigate whether there is a difference between the mean weight of male aquatic warblers and the mean weight of females.

Exercise 3.2 *Weights of fieldfares*

Use the data in Table 3.4 to investigate whether there is a difference between:

- the mean weight of adult male fieldfares and the mean weight of adult female fieldfares;
- the mean weight of first-year male fieldfares and the mean weight of first-year female fieldfares;
- the mean weight of first-year male fieldfares and the mean weight of adult male fieldfares.

4 Testing for a difference

To study this section, you will need access to your computer and the statistics software.



In Section 3, you saw how the two-sample z -test may be used to test for a difference between two population means, given two large samples of data from the populations. However, in no case did you have to do all the calculations: the sample means and sample standard deviations were provided. In practice, given two samples of data, you would have to calculate these statistics yourself. Alternatively, you can use a statistics software package to perform the test: the sample means and sample standard deviations will then be calculated automatically as part of the test. In this section, the use of OUStats to carry out a two-sample z -test is explained.

Refer to Computer Book D for the work in this section.



Summary of Section 4

This section has introduced the use of OUStats to do the calculations required to carry out a two-sample z -test. Large samples of data have been compared using boxplots, and the two-sample z -test has been used to test for a difference between two population means.

Summary of Chapter D4

In this chapter, two methods of comparing samples of data have been discussed. First, we reviewed the use of boxplots for comparing samples of data. Then the idea of a hypothesis test was introduced, in order to investigate whether there is a difference between the mean wing lengths of male and female meadow pipits. The test described was the two-sample z -test; this can be used only when both samples of data are large (at least 25). The importance of stating your hypotheses and conclusions clearly was stressed.

Learning outcomes

You have been working towards the following learning outcomes.

Terms to know and use

Median, lower quartile, upper quartile, range, interquartile range, boxplot, null and alternative hypotheses, standard error, estimated standard error, test statistic, significance level.

Symbols and notation to know and use

$Q1$ and $Q3$ for the lower quartile and the upper quartile;

H_0 for the null hypothesis of a statistical test and H_1 for the alternative hypothesis;

ESE for the estimated standard error of a sampling distribution.

Ideas to be aware of

- ◇ How boxplots drawn on a common axis may be used to compare two or more samples of data.
- ◇ That a hypothesis test consists of three stages: setting up the null and alternative hypotheses, calculating the test statistic, and reporting conclusions.
- ◇ How a 5% significance level is interpreted.
- ◇ How the difference between two sample means can be used to test whether or not there is a difference between the means of the populations from which the samples were drawn.

Features of OUStats

- ◇ Obtain boxplots to represent samples of data.
- ◇ Obtain the test statistic for a two-sample z -test given samples of data from two populations.

Solutions to Activities

Solution 1.2

All three methods have the advantage of being simple to use. The main disadvantage of the first method is that it does not distinguish between a participant who replaces the objects close to but not quite in their correct positions and one who replaces the objects very inaccurately on the grid. The second method also suffers from being a rather crude measure. The third method is more sophisticated than the other two. It has the advantage that it measures how far from its correct position each object is replaced, but has the disadvantage that it would be easy to forget that a low score indicates a good performance. However, none of the methods gives credit to a participant who places the objects in the correct pattern in relation to each other, but not in the correct positions – placing them all one square too high, for instance.

Solution 1.3

The city block scores of the group of 14 elderly people are written below in ascending order.

13 15 17 21 22 23 26
29 32 34 35 36 42 43

Since there is an even number of values, the median is the mean of the two middle values:

$$\text{median} = \frac{1}{2}(26 + 29) = 27.5$$

The lower quartile is the median of the values to the left of the median, and the upper quartile is the median of the values to the right of the median. This is illustrated in Figure S.1.

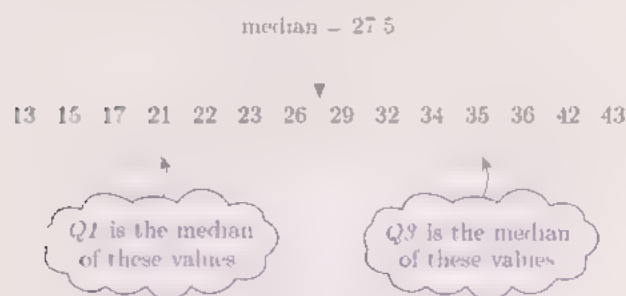


Figure S.1 Finding the median and the quartiles

From the figure, you can see that the lower quartile is 21 and the upper quartile is 35.

Solution 1.4

For the young group,

$$\text{range} = 36 - 4 = 32,$$

$$\text{interquartile range} = 21 - 6 = 15.$$

For the elderly group,

$$\text{range} = 43 - 13 = 30,$$

$$\text{interquartile range} = 35 - 21 = 14.$$

The values of the two measures of spread are roughly equal for the two groups, indicating that the spread of the city block scores is similar for the two groups.

Solution 1.5

- (a) The memorisation times of the 13 young people are shown in ascending order in Figure S.2.

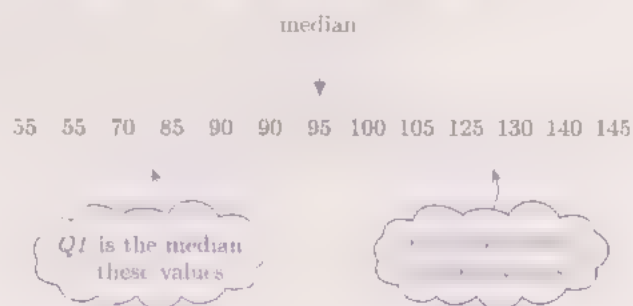


Figure S.2 Finding the median and the quartiles

For the young people, the median is 95.

The quartiles are

$$Q1 = \frac{1}{2}(70 + 85) = 77.5,$$

$$Q3 = \frac{1}{2}(125 + 130) = 127.5.$$

The memorisation times of the 14 elderly people are shown in ascending order in Figure S.3.

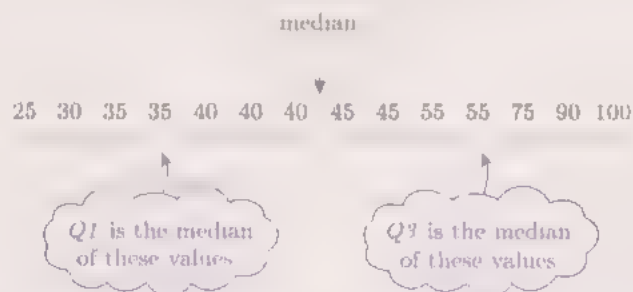


Figure S.3 Finding the median and the quartiles

For the elderly people, the median and quartiles are as follows:

$$\text{median} = \frac{1}{2}(40 + 45) = 42.5,$$

$$Q1 = 35, \quad Q3 = 55$$

- (b) Boxplots for the memorisation times in seconds of the young people and the elderly people are shown in Figure S.4.

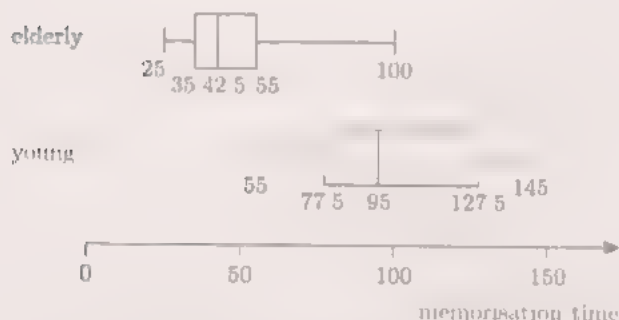


Figure S.4 Boxplots for the memorisation times

- (c) From the boxplots, it can be seen that the young people generally spent longer memorising the positions of the objects. All five key values on a boxplot – the minimum, the lower quartile, the median, the upper quartile and the maximum – are higher for the young group. In particular, notice that approximately three-quarters of the elderly people spent less time studying the positions of the objects than any of the young people.

It is also evident from the boxplots that the times spent by the elderly group were less widely spread than the times spent by the young group. For instance, the interquartile range is 50 seconds for the young group, but only 20 seconds for the elderly group. (Compare the lengths of the boxes.)

The boxplots suggest that the young people spent longer studying the positions of the objects than did the elderly people. You saw earlier that the young people also performed better on the test. So it is possible that the young people remembered the positions of the objects more accurately because they had spent longer memorising them. You will have the opportunity to investigate this in Section 2, using OU'Stats.

Solution 3.2

- (a) The estimated standard error is

$$ESE = \sqrt{\frac{s_M^2}{n_M} + \frac{s_F^2}{n_F}} = \sqrt{\frac{1.79^2}{31} + \frac{2.15^2}{27}} = 0.523986\dots \approx 0.52$$

- (b) The test statistic is

$$z = \frac{\bar{x}_M - \bar{x}_F}{ESE} = \frac{81.5 - 77.5}{0.523986\dots} \approx 7.63.$$

- (c) The test statistic is 'large': $z = 7.63 > 1.96$. So we reject the null hypothesis H_0 in favour of the alternative hypothesis H_1 .

Solution 3.3

- (a) The null and alternative hypotheses may be written as

$$H_0: \mu_{YM} = \mu_{AM},$$

$$H_1: \mu_{YM} \neq \mu_{AM},$$

where μ_{YM} is the mean wing length of the population of first-year male fieldfares and μ_{AM} is the mean wing length of the population of adult male fieldfares. (You may have used subscripts other than these for the two populations.)

The estimated standard error of the difference between two sample means is

$$ESE = \sqrt{\frac{s_{YM}^2}{n_{YM}} + \frac{s_{AM}^2}{n_{AM}}} = \sqrt{\frac{3.10^2}{131} + \frac{3.19^2}{80}} = 0.447839\dots,$$

and the test statistic is

$$z = \frac{\bar{x}_{YM} - \bar{x}_{AM}}{ESE} = \frac{150.0 - 151.9}{0.447839\dots} \approx -4.24$$

Since the test statistic is $z = -4.24 < -1.96$, we reject the null hypothesis at the 5% significance level in favour of the alternative hypothesis. We conclude that the mean wing length of first-year male fieldfares is not equal to the mean wing length of adult male fieldfares. The sample mean is greater for the adult males than for the first-year males, so this suggests that the mean wing length of adult males is greater than the mean wing length of first-year males.

(b) The null and alternative hypotheses may be written as

$$\begin{aligned} H_0 : \mu_{YF} &= \mu_{AF}, \\ H_1 : \mu_{YF} &\neq \mu_{AF}, \end{aligned}$$

where μ_{YF} is the mean wing length of the population of first-year female fieldfares and μ_{AF} is the mean wing length of the population of adult female fieldfares. (Again, you may have used subscripts other than these for the two populations.)

The estimated standard error is

$$ESE = \sqrt{\frac{s_{YF}^2}{n_{YF}} + \frac{s_{AF}^2}{n_{AF}}} = \sqrt{\frac{3.37^2}{288} + \frac{3.37^2}{139}} \\ = 0.365051\dots,$$

and the test statistic is

$$z = \frac{\bar{x}_{YF} - \bar{x}_{AF}}{ESE} = \frac{146.1 - 147.5}{0.365051} \approx -3.84$$

Since the test statistic is $z = -3.84 < -1.96$, we reject the null hypothesis at the 5% significance level in favour of the alternative hypothesis. We conclude that the mean wing length of first-year female fieldfares is not equal to the mean wing length of adult female fieldfares. The sample mean is greater for the adult females than for the first-year females, so this suggests that the mean wing length of adult females is greater than the mean wing length of first-year females.

Solution 3.4

The null and alternative hypotheses may be written as

$$\begin{aligned} H_0 : \mu_{YF} &= \mu_{AF}, \\ H_1 : \mu_{YF} &\neq \mu_{AF}, \end{aligned}$$

where μ_{YF} is now the mean weight of the population of first-year female fieldfares and μ_{AF} is the mean weight of the population of adult female fieldfares. (You may have used subscripts other than these for the two populations.)

The estimated standard error is

$$ESE = \sqrt{\frac{s_{YF}^2}{n_{YF}} + \frac{s_{AF}^2}{n_{AF}}} = \sqrt{\frac{9.52^2}{288} + \frac{9.06^2}{139}} \\ = 0.951429\dots,$$

and the test statistic is

$$\frac{\bar{x}_{YF} - \bar{x}_{AF}}{ESE} = \frac{108.0 - 108.7}{0.951429\dots} \approx -0.74.$$

Since $-1.96 < z < 1.96$, we cannot reject the null hypothesis at the 5% significance level. There is no evidence to suggest that the mean weight of first-year female fieldfares is different from the mean weight of adult female fieldfares.

Solution 3.5

You might add something like the following to the table you began in Chapter D3

Terminology	Notation	Standard deviation of...	Useful results
Standard error of the difference between two sample means	SE	sampling distribution of the difference between two sample means	$SE = \sqrt{\frac{\sigma_A^2}{n_A} + \frac{\sigma_B^2}{n_B}}$ where n_A, n_B are the sample sizes
	σ_A	population A	
	σ_B	population B	
Estimated standard error of the difference between two sample means	ESE		$ESE = \sqrt{\frac{s_A^2}{n_A} + \frac{s_B^2}{n_B}}$
	s_A	sample from A	
	s_B	sample from B	

In a two-sample z -test, the test statistic is

$$\frac{\bar{x}_A - \bar{x}_B}{ESE}$$

where \bar{x}_A and \bar{x}_B are the means of the samples from populations A and B.

Solutions to Exercises

Solution 1.1

- (a) The gross weekly earnings of the 9 female police officers are shown in ascending order in Figure S.5.

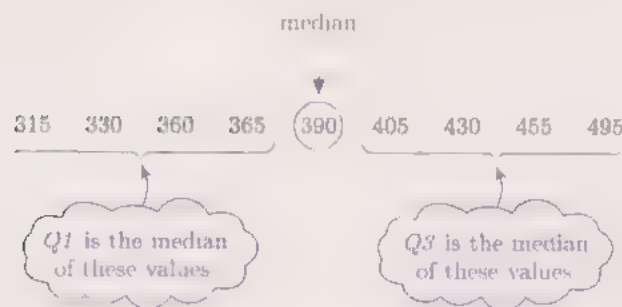


Figure S.5 Finding the median and the quartiles

For the female police officers, the median is 390
The quartiles are

$$Q1 = \frac{1}{2}(330 + 360) = 345,$$
$$Q3 = \frac{1}{2}(430 + 455) = 442.5.$$

The gross weekly earnings of the 10 male police officers are shown in ascending order in Figure S.6.

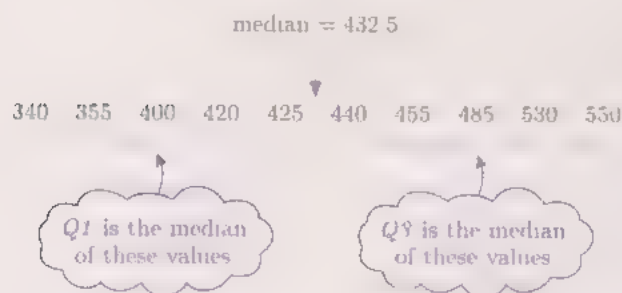


Figure S.6 Finding the median and the quartiles

For the male police officers, the median and quartiles are as follows:

$$\text{median} = \frac{1}{2}(425 + 440) = 432.5,$$
$$Q1 = 400, \quad Q3 = 485.$$

Boxplots for the gross weekly earnings in pounds of the male and female police officers are shown in Figure S.7.

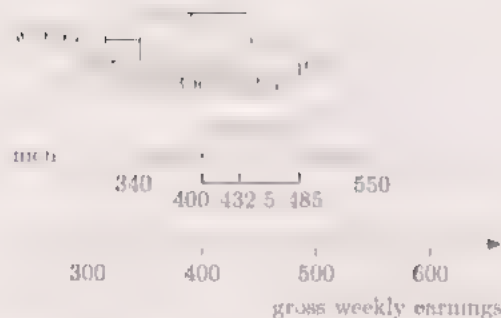


Figure S.7 Boxplots for the earnings of male and female police officers

- (c) From the boxplots, it can be seen that the earnings of the men were generally a little higher than the earnings of the women. All five key values on a boxplot – the minimum, the lower quartile, the median, the upper quartile and the maximum – are higher for the men than for the women, though not by very much.
- (d) For the women:

$$\text{range} = 495 - 315 = 180,$$
$$\text{interquartile range} = 442.5 - 345 = 97.5.$$

For the men:

$$\text{range} = 550 - 340 = 210,$$
$$\text{interquartile range} = 485 - 400 = 85.$$

The values of the two measures of spread are similar for the two groups, indicating that the spread of earnings is similar for the men and the women.

Solution 1.2

- (a) The gross hourly earnings of the 8 female chefs and cooks are shown in ascending order in Figure S.8.

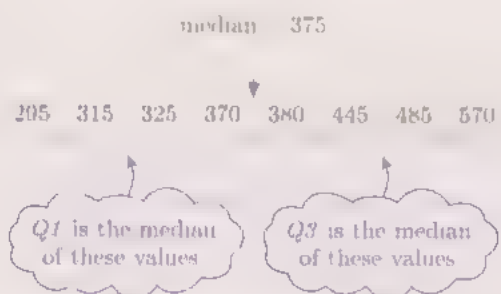


Figure S.8 Finding the median and the quartiles

For the female chefs and cooks, the median is $\frac{1}{2}(370 + 380) = 375$.

The quartiles are

$$Q1 = \frac{1}{2}(315 + 325) = 320,$$

$$Q3 = \frac{1}{2}(445 + 485) = 465.$$

The gross hourly earnings of the 11 male chefs and cooks are shown in ascending order in Figure S.9.

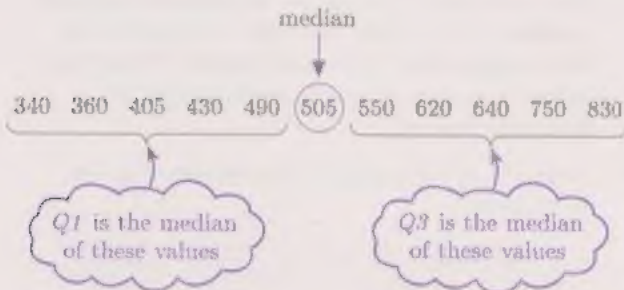


Figure S.9 Finding the median and the quartiles

For the male chefs and cooks, the median and quartiles are as follows:

$$\text{median} = 505,$$

$$Q1 = 405, \quad Q3 = 640.$$

- (b) Boxplots for the gross hourly earnings in pence of the male and female chefs and cooks are shown in Figure S.10.

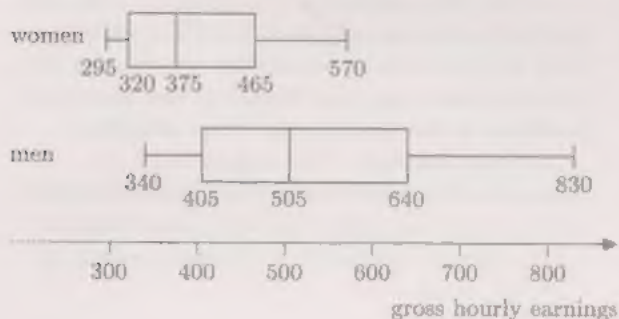


Figure S.10 Boxplots for the earnings of male and female chefs and cooks

- (c) From the boxplots, it can be seen that the earnings of the men were generally higher than the earnings of the women. All five key values on a boxplot – the minimum, the lower quartile, the median, the upper quartile and the maximum – are higher for the men than for the women. In particular, the lower quartile of the men's earnings is greater than the median earnings of the women; and all the men earned more than the lower quartile of the women's earnings (indicating that more than a quarter of the women earned less than any of the men).

- (d) For the women:

$$\text{range} = 570 - 295 = 275;$$

$$\text{interquartile range} = 465 - 320 = 145.$$

For the men:

$$\text{range} = 830 - 340 = 490;$$

$$\text{interquartile range} = 640 - 405 = 235.$$

The spread of the men's earnings is much greater than the spread of the women's earnings. Both the range and the interquartile range are greater for the men than for the women.

Solution 3.1

The null and alternative hypotheses may be written as

$$H_0: \mu_M = \mu_F,$$

$$H_1: \mu_M \neq \mu_F,$$

where μ_M is the mean weight of the population of male aquatic warblers and μ_F is the mean weight of the population of female aquatic warblers. (You may have used subscripts other than these for the two populations.)

The estimated standard error is

$$\begin{aligned} ESE &= \sqrt{\frac{s_M^2}{n_M} + \frac{s_F^2}{n_F}} = \sqrt{\frac{0.73^2}{66} + \frac{0.87^2}{83}} \\ &= 0.131\,124\dots, \end{aligned}$$

and the test statistic is

$$z = \frac{\bar{x}_M - \bar{x}_F}{ESE} = \frac{12.6 - 12.1}{0.131\,124\dots} \approx 3.81.$$

Since the test statistic is $z \approx 3.81 > 1.96$, we reject the null hypothesis at the 5% significance level in favour of the alternative hypothesis. We conclude that the mean weight of male aquatic warblers is not equal to the mean weight of female aquatic warblers. The sample mean is greater for the males than for the females, so this suggests that the mean weight of male aquatic warblers is greater than the mean weight of female aquatic warblers.

Solution 3.2

- (a) The null and alternative hypotheses may be written as

$$H_0: \mu_{AM} = \mu_{AF},$$

$$H_1: \mu_{AM} \neq \mu_{AF},$$

where μ_{AM} is the mean weight of the population of adult male fieldfares and μ_{AF} is the mean weight of the population of adult female fieldfares. (You may have used subscripts other than these for the two populations.)

The estimated standard error is

$$\begin{aligned} ESE &= \sqrt{\frac{s_{AM}^2}{n_{AM}} + \frac{s_{AF}^2}{n_{AF}}} = \sqrt{\frac{10.91^2}{93} + \frac{9.06^2}{139}} \\ &= 1.367626\dots, \end{aligned}$$

and the test statistic is

$$z = \frac{\bar{x}_{AM} - \bar{x}_{AF}}{ESE} = \frac{114.9 - 108.7}{1.367626\dots} \simeq 4.53.$$

Since the test statistic is $z = 4.53 > 1.96$, we reject the null hypothesis at the 5% significance level in favour of the alternative hypothesis. We conclude that the mean weight of adult male fieldfares is not equal to the mean weight of adult female fieldfares. The sample mean is greater for the adult males than for the adult females, so this suggests that the mean weight of adult males is greater than the mean weight of adult females.

- (b) The null and alternative hypotheses may be written as

$$H_0: \mu_{YM} = \mu_{YF},$$

$$H_1: \mu_{YM} \neq \mu_{YF},$$

where μ_{YM} is the mean weight of the population of first-year male fieldfares and μ_{YF} is the mean weight of the population of first-year female fieldfares.

The estimated standard error is

$$\begin{aligned} ESE &= \sqrt{\frac{s_{YM}^2}{n_{YM}} + \frac{s_{YF}^2}{n_{YF}}} = \sqrt{\frac{8.62^2}{144} + \frac{9.52^2}{288}} \\ &= 0.911422\dots, \end{aligned}$$

and the test statistic is

$$z = \frac{\bar{x}_{YM} - \bar{x}_{YF}}{ESE} = \frac{111.6 - 108.0}{0.911422\dots} \simeq 3.95.$$

Since the test statistic is $z = 3.95 > 1.96$, we reject the null hypothesis at the 5% significance level in favour of the alternative hypothesis. We conclude that the mean weight of first-year male fieldfares is not equal to the mean weight of first-year female fieldfares. The sample mean is greater for the males than for the females, so this suggests that the mean weight of first-year males is greater than the mean weight of first-year females.

- (c) The null and alternative hypotheses may be written as

$$H_0: \mu_{YM} = \mu_{AM},$$

$$H_1: \mu_{YM} \neq \mu_{AM},$$

where μ_{YM} is the mean weight of the population of first-year male fieldfares and μ_{AM} is the mean weight of the population of adult male fieldfares.

The estimated standard error is

$$\begin{aligned} ESE &= \sqrt{\frac{s_{YM}^2}{n_{YM}} + \frac{s_{AM}^2}{n_{AM}}} = \sqrt{\frac{8.62^2}{144} + \frac{10.91^2}{93}} \\ &= 1.340102\dots, \end{aligned}$$

and the test statistic is

$$z = \frac{\bar{x}_{YM} - \bar{x}_{AM}}{ESE} = \frac{111.6 - 114.9}{1.340102\dots} \simeq -2.46.$$

Since the test statistic is $z = -2.46 < -1.96$, we reject the null hypothesis at the 5% significance level in favour of the alternative hypothesis. We conclude that the mean weight of first-year male fieldfares is not equal to the mean weight of adult male fieldfares. The sample mean is greater for the adult males than for the first-year males, so this suggests that the mean weight of adult males is greater than the mean weight of first-year males.



Using Mathematics

BLOCK A **MODELLING WITH MATHEMATICS**

CHAPTER A1 *Modelling physical processes*

CHAPTER A2 *Modelling growth*

CHAPTER A3 *Representing circles*

CHAPTER A4 *Modelling with functions*

COMPUTER BOOK A

BLOCK B **DISCRETE MODELS**

CHAPTER B1 *Functions and calculations*

CHAPTER B2 *Modelling with sequences*

CHAPTER B3 *Modelling with matrices*

COMPUTER BOOK B

BLOCK C **CONTINUOUS MODELS**

CHAPTER C1 *Differentiation and modelling*

CHAPTER C2 *Integration and modelling*

CHAPTER C3 *Choosing a function for a model*

COMPUTER BOOK C

BLOCK D **MODELLING UNCERTAINTY**

CHAPTER D1 *Chance*

CHAPTER D2 *Modelling variation*

CHAPTER D3 *Estimating*

CHAPTER D4 *Comparing*

CHAPTER D5 *Looking for relationships*

COMPUTER BOOK D